

Indécidabilité de la satisfiabilité d'une requête

LEÇONS : 914 ; 932

RÉFÉRENCES : ABITEBOUL–HULL–VIANU, *Foundations of Database* (p.123) [?] et LESESVRE–MONTAGNON–LE BARBENCHON–PIERRON, *131 développements pour l'oral* (p. 817) [?]

Notations :

- pour $u_{i_1} \dots u_{i_n}$, on note $u^{(k)}$ pour $u_{i_1} \dots u_{i_k}$ avec $k \in \llbracket 1, n \rrbracket$
- pour $v_{i_1} \dots v_{i_n}$, on note $v^{(k)}$ pour $v_{i_1} \dots v_{i_k}$ avec $k \in \llbracket 1, n \rrbracket$

Introduction :

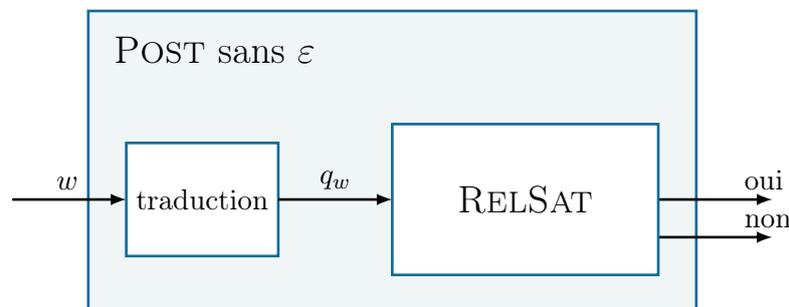
On veut montrer que le problème qui est de savoir si une requête du calcul relationnel est satisfiable est un problème indécidable. On va pour le prouver montrer une réduction du problème POST à RELSAT.

Théorème 1. Le problème

RELSAT	{	entrée : une requête q du calcul relationnel
	}	sortie : oui si il existe une base de donnée \mathcal{I} telle que $q(\mathcal{I}) \neq \emptyset$, non sinon

est indécidable.

Démonstration. On va réduire le problème POST sans ε ¹ à RELSAT.



Soit w une instance de POST, c'est à dire un nombre fini de tuiles $\begin{matrix} u_1 \\ v_1 \end{matrix}, \dots, \begin{matrix} u_N \\ v_N \end{matrix}$ sur l'alphabet $\Sigma = \{a, b\}$.

On prend le schéma de base de données suivant

$$\mathcal{R} = \{\text{ENC}(pos1, pos2, mot, u, v), \text{SYNC}(debut_u, debut_v)\}$$

Je note ENC pour ENCODER et SYNC pour SYNCHRONISER

Les colonnes $pos1$ et $pos2$ servent à mettre un ordre sur les tuples. Pour une instance positive de POST, on écrira le mot qui figure en haut et en bas d'une concaténation de tuiles valide dans la colonne mot . La colonne u (respectivement v) sert à donner l'indice de la tuile du haut (resp. du bas) à laquelle correspond la lettre qui est dans la colonne mot . (voir les exemples de remplissage du tableau plus bas)

1.

POST sans ε	{	entrée : des tuiles du problème de POST sans tuile contenant des ε
	}	sortie : oui si on peut trouver une suite de tuiles ayant le même mot en haut et en bas, non sinon

(pour la réduction du problème de Post original voir l'appendice [ici](#)).

On pose la requête suivante

$$q_w = \{\langle \rangle, \varphi_w\}$$

où φ_w est la conjonction des 5 formules suivantes :

- φ_{ENC} qui traduit le fait que
 - les éléments de $pos1$ sont disjoints
 - les colonnes $pos1, pos2$ forment un cycle de longueur > 1
 - $pos1$ contient le symbole $\$$

Cette formule sert à forcer les tuples de la relation ENC à être dans un certain ordre² et donc de retrouver cet ordre même si une relation est un ensemble de tuples (donc sans ordre a priori). C'est pour ça que l'on a besoin de deux colonnes. La première colonne est une clé pour retrouver le début de chaque tuile (d'où la nécessité d'avoir des éléments distincts). Le symbole $\$$ est un symbole spécial pour marquer le début.

- φ_{SYNC} qui traduit le fait que les éléments de $debut_u$ sont disjoints et les éléments de $debut_v$ sont disjoints

On a les indices de début de chaque tuile, donc il faut que les éléments soient distincts.

- φ_u qui traduit le fait que pour tout $x \in debut_u$, il existe $i \in \llbracket 1, N \rrbracket$ avec $u_i = z_1 \dots z_k$ et x_1, \dots, x_{k+1} avec $x = x_1$ tels que

$$\langle x_1, x_2, z_1, i, * \rangle \in \text{ENC}$$

$$\langle x_2, x_3, z_2, i, * \rangle \in \text{ENC}$$

$$\vdots$$

$$\langle x_k, x_{k+1}, z_k, i, * \rangle \in \text{ENC}$$

avec $x_1, x_{k+1} \in debut_u$ et $x_2, \dots, x_k \notin debut_u$
(où $*$ peut être choisi indifféremment)

Cette formule permet de dire que pour chaque indice de début de tuile dans SYNC, on a

une tuile $\begin{array}{|c|} \hline u_i \\ \hline v_i \\ \hline \end{array}$ et que (dans l'ordre imposé par la formule φ_{ENC}) on a bien une succession de

tuples qui correspond au mot u_i de la tuile $\begin{array}{|c|} \hline u_i \\ \hline v_i \\ \hline \end{array}$ puis que les lettres du mot correspondent bien à u_i entre les indices $x = x_1$ et x_k . De plus, on oblige à ne pas avoir de début de tuile pendant les indices intermédiaires, et que l'indice suivant (x_{k+1}) soit un début de tuile.

- φ_v similaire à φ_u

φ_v qui traduit le fait que pour tout $x \in debut_v$, il existe $i \in \llbracket 1, N \rrbracket$ avec $v_i = z_1 \dots z_k$ et x_1, \dots, x_{k+1} avec $x = x_1$ tels que

$$\langle x_1, x_2, z_1, *, i \rangle \in \text{ENC}$$

$$\langle x_2, x_3, z_2, *, i \rangle \in \text{ENC}$$

$$\vdots$$

$$\langle x_k, x_{k+1}, z_k, *, i \rangle \in \text{ENC}$$

avec $x_1, x_{k+1} \in debut_v$ et $x_2, \dots, x_k \notin debut_v$
(où $*$ peut être choisi indifféremment)

2. l'ordre que l'on donne dans les tables que l'on dessine dans le développement, c'est en fait l'ordre dans lequel on lit le mot qui est en haut et en bas d'une concaténation valide de tuiles du problème POST

Cette formule permet de dire que pour chaque indice de début de tuile dans SYNC, on a une tuile $\begin{array}{|c|} \hline u_i \\ \hline v_i \\ \hline \end{array}$ et que (dans l'ordre imposé par la formule φ_{ENC}) on a bien une succession de tuples qui correspond au mot v_i de la tuile $\begin{array}{|c|} \hline u_i \\ \hline v_i \\ \hline \end{array}$ puis que les lettres du mot correspondent bien à v_i entre les indices $x = x_1$ et x_k . De plus, on oblige à ne pas avoir de début de tuile pendant les indices intermédiaires, et que l'indice suivant (x_{k+1}) soit un début de tuile.

— $\varphi_{\text{SYNC}_u_v}$ définie par

SYNC(\$, \$)

$$\wedge \forall x \forall y \left(\text{SYNC}(x, y) \rightarrow \bigvee_{i=1}^N \exists t_1 t_2 t_3 s_1 s_2 s_3 \text{ENC}(x, t_1, t_2, i, t_3) \wedge \text{ENC}(y, s_1, s_2 s_3, i) \right)$$

$$\wedge \forall x \forall y \left(\text{SYNC}(x, y) \wedge \neg(x = \max(\text{debut}_u)) \wedge \neg(y = \max(\text{debut}_v)) \rightarrow \text{SYNC}(x', y') \right)$$

$$\text{avec } x' = \min_{a \in \text{debut}_u} a > x \text{ et } y' = \min_{a \in \text{debut}_v} a > y$$

Cette formule force à commencer par une première tuile (par le fait que SYNC(\$, \$)), puis à synchroniser le reste des tuiles, c'est-à-dire que après la tuile $\begin{array}{|c|} \hline u_{i_1} \\ \hline v_{i_1} \\ \hline \end{array}$, on accède à la tuile suivante (3^{ème} sous-formule), et que c'est une tuile licite (car il faut que ce soit une tuile $\begin{array}{|c|} \hline u_i \\ \hline v_i \\ \hline \end{array}$ où l'indice du haut est le même que l'indice du bas et non une tuile $\begin{array}{|c|} \hline u_3 \\ \hline v_4 \\ \hline \end{array}$ qui ne serait pas licite) (par la 2^{ème} sous-formule).

Il faut montrer que

$$\boxed{q_w \text{ est satisfiable SSI } w \text{ est une instance positive de POST}}$$

\Rightarrow Soit une instance \mathcal{I} telle que $\mathcal{I} \models q_w$.

On note $r = |\text{SYNC}|$ (ce sera le nombre de tuiles)

On veut montrer par récurrence sur $k \in \llbracket 1, r \rrbracket$ que

$$HR_k : \begin{cases} u^{(k)} = \text{mot}[1 \dots |u^{(k)}|] \\ v^{(k)} = \text{mot}[1 \dots |v^{(k)}|] \end{cases}$$

On prouve cela grâce aux formules de φ_w qui est satisfiable³.

Ainsi, on peut conclure que $u^{(r)} = v^{(r)}$ et donc que w est une instance positive de POST.

\Leftarrow w est une instance positive de POST, donc il existe i_1, \dots, i_r tels que

$$m = u_{i_1} \dots u_{i_r} = v_{i_1} \dots v_{i_r}$$

On construit alors la table ENC et SYNC comme ci-dessous

3. les arguments sont en fait les explications en gris sous chaque formule

ENC	<i>pos1</i>	<i>pos2</i>	<i>mot</i>	<i>u</i>	<i>v</i>
	\$	1	m_1	i_1	i_1
	1	2	m_2	i_1	i_1
	\vdots	\vdots	\vdots	i_1	i_1
	$ v^{(1)} $	$ v^{(1)} + 1$	$m_{ v^{(1)} }$	i_1	i_2
	\vdots	\vdots	\vdots	i_1	\vdots
	\vdots	\vdots	\vdots	i_1	\vdots
	$ u^{(1)} $	$ u^{(1)} + 1$	$m_{ u^{(1)} }$	i_2	\vdots
	\vdots	\vdots	\vdots	\vdots	\vdots
	\vdots	\vdots	\vdots	\vdots	\vdots
	$ m - 2$	$ m - 1$	\vdots	\vdots	\vdots
	$ m - 1$	\$	$m_{ m }$	i_r	i_r

SYNC	<i>debut_u</i>	<i>debut_v</i>
	\$	\$
	$ u^{(1)} $	$ v^{(1)} $
	$ u^{(2)} $	$ v^{(2)} $
	\vdots	\vdots
	\vdots	\vdots
	$ u^{(r-1)} $	$ v^{(r-1)} $

Cette instance de base de données \mathcal{I} (sur le schéma de base de données \mathcal{R}) satisfait la requête q_w ⁴. Il suffit de vérifier qu'elle vérifie chaque formule que l'on a énoncé. □

Remarques :

Il existe d'autres problèmes indécidables sur les bases de données comme :

<i>RelEqu</i>	{	entrée : deux requêtes q et q' du calcul relationnel
sortie : oui si q et q' sont sémantiquement équivalentes, non sinon		
<i>RelSure</i>	{	entrée : une requête q du calcul relationnel
sortie : oui si q est sûre, non sinon		

Pour montrer l'indécidabilité de RELEQU, il suffit de réduire co-RELSAT à RELEQU. En effet, une requête est insatisfiable si et seulement si elle est équivalente à la requête vide.

Astuces de l'agregatif :

- Ce développement est très long, je ne pense pas qu'il soit faisable d'écrire toutes les formules et de prouver la réduction tout en faisant comprendre au jury ce que l'on est en train de faire. J'ai fait le choix d'essayer de faire comprendre au jury la réduction en donnant le sens des formules.
- J'écris la base de données à droite du tableau dès que j'introduis ENC et SYNC en expliquant comment les remplir, pour que le sens des formules soit compréhensible quand je les énonce ensuite. Je dis que le but des formules est que la seule façon de remplir la base de données soit celle que j'ai présenté.
- Lors du développement je dis juste POST, en effet ce n'est qu'une variante de POST.
- Je mets les formules ci-dessous, il peut être intéressant de donner une formule d'unicité sur la formule φ_{ENC} pour convaincre le jury que l'on maîtrise aussi l'écriture rigoureuse des formules.
- On a prouvé ici que ce problème n'est pas dans R , mais il est quand même dans RE . En effet, on peut dénombrer les bases de données qui sont des ensembles finis de tuples sur un domaine dénombrable.

4. i.e. $q_w(\mathcal{I}) \neq \emptyset$

— φ_{ENC}

$$\begin{aligned}
& \exists t \left(\psi(\$, t) \wedge \neg(t = \$) \right) \wedge \exists t \left(\psi(t, \$) \wedge \neg(t = \$) \right) && \text{existence \$} \\
& \wedge \forall x \forall y \exists z \left(\psi(x, y) \rightarrow \psi(y, z) \right) && \text{cycle} \\
& \wedge \forall x \forall y \exists z \left(\psi(y, x) \rightarrow \psi(z, y) \right) && \text{cycle} \\
& \wedge \forall x \forall y \forall z \left((\psi(x, y) \wedge \psi(x, z)) \rightarrow (y = z) \right) && \text{unicité} \\
& \wedge \forall x \forall y \forall z \left((\psi(y, x) \wedge \psi(z, x)) \rightarrow (y = z) \right) && \text{unicité} \\
& \wedge \forall x t_1 t_2 t_3 t_4 s_1 s_2 s_3 s_4 \left((\text{ENC}(x, t_1, t_2, t_3, t_4) \right. \\
& \quad \left. \wedge \text{ENC}(x, s_1, s_2, s_3, s_4)) \rightarrow (t_1 = s_1 \right. \\
& \quad \quad \wedge t_2 = s_2 \\
& \quad \quad \wedge t_3 = s_3 \\
& \quad \quad \left. \wedge t_4 = s_4) \right) && \text{unicité } \textit{pos1}
\end{aligned}$$

avec $\psi(x, y) := \exists z \exists t \exists s \text{ENC}(x, y, z, t, s)$ — φ_{SYNC}

$$\begin{aligned}
& \forall x \forall y \forall z \left((\text{SYNC}(x, y) \wedge \text{SYNC}(x, z)) \rightarrow y = z \right) && \text{unicité } \textit{debut_u} \\
& \wedge \forall x \forall y \forall z \left((\text{SYNC}(x, y) \wedge \text{SYNC}(z, y)) \rightarrow x = z \right) && \text{unicité } \textit{debut_v}
\end{aligned}$$

— φ_u

$$\begin{aligned}
\forall x \left(\exists y \text{SYNC}(x, y) \rightarrow \left(\bigvee_{i=1}^N \left(\exists x_1 x_2 \dots x_{|u_i|} x_{|u_i|+1} j_1 \dots j_{|u_i|} \right. \right. \right. \\
\quad \left. \left. \left(\bigwedge_{k=1}^{|u_i|} \text{ENC}(x_k, x_{k+1}, u_i[k], i, j_k) \right) \right. \right. \\
\quad \wedge (x = x_1) \wedge (\exists z \text{SYNC}(x_{|u_i|+1}, z)) \\
\quad \left. \left. \left. \wedge \bigwedge_{k=2}^{|u_i|} \neg \exists y \text{SYNC}(x_k, y) \right) \right) \right)
\end{aligned}$$

— φ_v similaire

— $\varphi_{\text{SYNC_u_v}}$ il faut coder $x < y$, pour cela il faut borner le nombre d'intermédiaire entre x et y , on n'a a priori pas de borne, mais on sait que l'indice de début de tuile suivante est forcément à moins de "taille maximum des mots sur les tuiles" de x , donc on peut faire un disjonction sur le nombre d'intermédiaires qui sera donc borné puisqu'on regarde toujours la différence entre x et y sur des tuiles qui se suivent. Pour faire l'égalité avec le maximum, il faut regarder si dans le tableau ENC, le dernier caractère de la tuile possède une *pos2* égale à $\$$.

$$\begin{aligned}
x < y &:= \neg(x = y) \wedge \bigvee_{k=1}^{\text{borne}} \exists z_1 z_2 \dots z_k \left(\psi(x, z_1) \wedge \psi(z_1, z_2) \wedge \dots \wedge \psi(z_k, y) \right) \\
x = \max(\textit{debut_u}) &:= \bigvee_{k=0}^{|u_x|} \exists z_1 \dots z_k \left(\psi(x, z_1) \wedge \psi(z_1, z_2) \wedge \dots \wedge \psi(z_k, \$) \right)
\end{aligned}$$

(Je ne suis pas sûr de ces dernières formules, il faudrait écrire vraiment rigoureusement tout cela pour être vraiment propre, mais je pense que cette idée marche⁵).

Je mets un exemple concret de tables pour $m = aabbbab = u_{i_1}u_{i_2}u_{i_3} = v_{i_1}v_{i_2}v_{i_3}$ avec les tuiles

$$i_1 \begin{array}{|c|} \hline aab \\ \hline a \\ \hline \end{array}; i_2 \begin{array}{|c|} \hline b \\ \hline abbb \\ \hline \end{array}; i_3 \begin{array}{|c|} \hline bab \\ \hline ab \\ \hline \end{array}$$

ENC

<i>pos1</i>	<i>pos2</i>	<i>mot</i>	<i>u</i>	<i>v</i>
\$	1	<i>a</i>	i_1	i_1
1	2	<i>a</i>	i_1	i_2
2	3	<i>b</i>	i_1	i_2
3	4	<i>b</i>	i_2	i_2
4	5	<i>b</i>	i_3	i_2
5	6	<i>a</i>	i_3	i_3
6	\$	<i>b</i>	i_3	i_3

SYNC

<i>debut_u</i>	<i>debut_v</i>
\$	\$
3	1
4	5

Questions possibles :

- Montrer que le problème de POST sans ε est indécidable.
- Écrire explicitement la récurrence de la question **c)** (*i*).

5. envoie moi un [mail](#), si tu réussis à écrire ça proprement